

Aim: Find the Outlier from the given data set using trimming and capping methods.

```
In [18]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib.pyplot as plt

%matplotlib inline
```

```
In [19]: df = pd.read_csv('placement.csv')
df
```

```
Out[19]:
```

	cgpa	placement_exam_marks	placed
0	7.19	26.0	1
1	7.46	38.0	1
2	7.54	40.0	1
3	6.42	8.0	1
4	7.23	17.0	0
...
995	8.87	44.0	1
996	9.12	65.0	1
997	4.89	34.0	0
998	8.62	46.0	1
999	4.90	10.0	1

1000 rows × 3 columns

```
In [20]: df.head()
```

```
Out[20]:
```

	cgpa	placement_exam_marks	placed
0	7.19	26.0	1
1	7.46	38.0	1
2	7.54	40.0	1
3	6.42	8.0	1
4	7.23	17.0	0

```
In [21]: df.tail()
```

```
Out[21]:
```

	cgpa	placement_exam_marks	placed
995	8.87	44.0	1
996	9.12	65.0	1
997	4.89	34.0	0
998	8.62	46.0	1
999	4.90	10.0	1

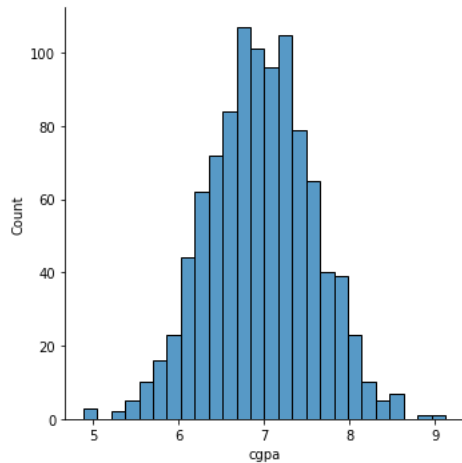
```
In [22]: %matplotlib.notebook
```

UsageError: Line magic function `%matplotlib.notebook` not found.

In [23]:

```
sns.displot(df['cgpa'])
```

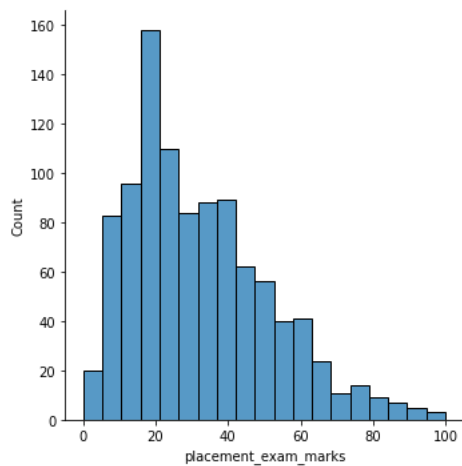
Out[23]: <seaborn.axisgrid.FacetGrid at 0x163a4803eb0>



In [27]:

```
sns.displot(df['placement_exam_marks'])
```

Out[27]: <seaborn.axisgrid.FacetGrid at 0x163a497f910>



In [28]:

```
df['placement_exam_marks'].describe()
```

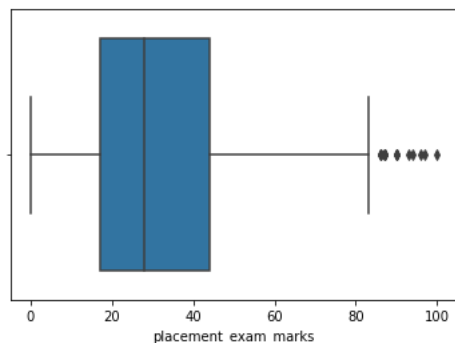
```
Out[28]: count    1000.000000
mean       32.225000
std        19.130822
min         0.000000
25%        17.000000
50%        28.000000
75%        44.000000
max        100.000000
Name: placement_exam_marks, dtype: float64
```

```
In [29]: sns.boxplot(df['placement_exam_marks'])
```

C:\Users\User38\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```

```
Out[29]: <AxesSubplot:xlabel='placement_exam_marks'>
```



```
In [30]: # finding highest boundaries values
print('Highest Boundary value of Cgpa', df['cgpa'].mean() + 3*df['cgpa'].std())
```

```
Highest Boundary value of Cgpa 8.808933625397177
```

```
In [31]: # Finding Lowest boundaries value
print('Lowest Boundary value of Cgpa', df['cgpa'].mean() - 3*df['cgpa'].std())
```

```
Lowest Boundary value of Cgpa 5.113546374602842
```

```
In [32]: # finding outliers
df[(df['cgpa']>8.80) | (df['cgpa']<5.11)]
```

```
Out[32]:
```

	cgpa	placement_exam_marks	placed
485	4.92	44.0	1
995	8.87	44.0	1
996	9.12	65.0	1
997	4.89	34.0	0
999	4.90	10.0	1

Trimming

```
In [33]: df.shape
```

```
Out[33]: (1000, 3)
```

```
In [34]: new_df = df[(df['cgpa']<8.80) & (df['cgpa']>5.11)]
new_df
```

```
Out[34]:
```

	cgpa	placement_exam_marks	placed
0	7.19	26.0	1
1	7.46	38.0	1
2	7.54	40.0	1
3	6.42	8.0	1
4	7.23	17.0	0
...
991	7.04	57.0	0
992	6.26	12.0	0
993	6.73	21.0	1
994	6.48	63.0	0
998	8.62	46.0	1

```
995 rows × 3 columns
```

In [35]: `new_df.shape`

Out[35]: (995, 3)

Z Score

In [36]: `df['cgpa_score'] = (df['cgpa'] - df['cgpa'].mean())/df['cgpa'].std()
df`

Out[36]:

	cgpa	placement_exam_marks	placed	cgpa_score
0	7.19	26.0	1	0.371425
1	7.46	38.0	1	0.809810
2	7.54	40.0	1	0.939701
3	6.42	8.0	1	-0.878782
4	7.23	17.0	0	0.436371
...
995	8.87	44.0	1	3.099150
996	9.12	65.0	1	3.505062
997	4.89	34.0	0	-3.362960
998	8.62	46.0	1	2.693239
999	4.90	10.0	1	-3.346724

1000 rows × 4 columns

In [37]: `df.describe()`

Out[37]:

	cgpa	placement_exam_marks	placed	cgpa_score
count	1000.000000	1000.000000	1000.000000	1.000000e+03
mean	6.961240	32.225000	0.489000	-1.600275e-14
std	0.615898	19.130822	0.500129	1.000000e+00
min	4.890000	0.000000	0.000000	-3.362960e+00
25%	6.550000	17.000000	0.000000	-6.677081e-01
50%	6.960000	28.000000	0.000000	-2.013321e-03
75%	7.370000	44.000000	1.000000	6.636815e-01
max	9.120000	100.000000	1.000000	3.505062e+00

In [38]: `df['cgpa_score'].describe()`

Out[38]: count 1.000000e+03
mean -1.600275e-14
std 1.000000e+00
min -3.362960e+00
25% -6.677081e-01
50% -2.013321e-03
75% 6.636815e-01
max 3.505062e+00
Name: cgpa_score, dtype: float64

In [39]: `df[df['cgpa_score']>3]`

Out[39]:

	cgpa	placement_exam_marks	placed	cgpa_score
995	8.87	44.0	1	3.099150
996	9.12	65.0	1	3.505062

In [40]: `df[df['cgpa_score']< -3]`

Out[40]:

	cgpa	placement_exam_marks	placed	cgpa_score
485	4.92	44.0	1	-3.314251
997	4.89	34.0	0	-3.362960
999	4.90	10.0	1	-3.346724

```
In [41]: new_df = df[(df['cgpa_score']<3) & (df['cgpa_score']>-3)]
new_df.shape
```

```
Out[41]: (995, 4)
```

Capping

```
In [42]: upper_limit = df['cgpa'].mean() + 3*df['cgpa'].std()
lower_limit = df['cgpa'].mean() - 3*df['cgpa'].std()
lower_limit
```

```
Out[42]: 5.113546374602842
```

```
In [43]: df['cgpa_cap'] = np.where(
    df['cgpa']>upper_limit,
    upper_limit,
    np.where(
        df['cgpa']<lower_limit,
        lower_limit,df['cgpa']

    )

)
```

```
In [44]: df.describe()
```

```
Out[44]:
```

	cgpa	placement_exam_marks	placed	cgpa_score	cgpa_cap
count	1000.000000	1000.000000	1000.000000	1.000000e+03	1000.000000
mean	6.961240	32.225000	0.489000	-1.600275e-14	6.961499
std	0.615898	19.130822	0.500129	1.000000e+00	0.612688
min	4.890000	0.000000	0.000000	-3.362960e+00	5.113546
25%	6.550000	17.000000	0.000000	-6.677081e-01	6.550000
50%	6.960000	28.000000	0.000000	-2.013321e-03	6.960000
75%	7.370000	44.000000	1.000000	6.636815e-01	7.370000
max	9.120000	100.000000	1.000000	3.505062e+00	8.808934

```
In [45]: df['placement_exam_marks'].skew()
```

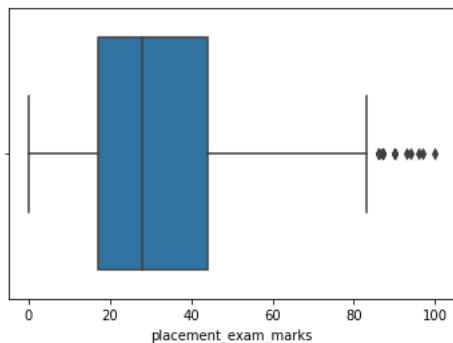
```
Out[45]: 0.8356419499466834
```

```
In [46]: sns.boxplot(df['placement_exam_marks'])
```

C:\Users\User38\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
warnings.warn(
```

```
Out[46]: <AxesSubplot:xlabel='placement_exam_marks'>
```



```
In [47]: q1,q2,q3=df['placement_exam_marks'].describe()[["25%", "50%", "75%"]]
```

```
In [49]:
```

```
q1
```

```
Out[49]: 17.0
```

```
In [50]: q2
```

```
Out[50]: 28.0
```

```
In [51]: q3
```

```
Out[51]: 44.0
```

```
In [52]: iqr = q3-q1
```

```
In [53]: iqr
```

```
Out[53]: 27.0
```

```
In [54]: upper_limit =q3 +1.5*iqr  
upper_limit
```

```
Out[54]: 84.5
```

```
In [55]:  
lower_limit =q1 -1.5*iqr  
lower_limit
```

```
Out[55]: -23.5
```

```
In [56]: df[df['placement_exam_marks'] > upper_limit].shape
```

```
Out[56]: (15, 5)
```

```
In [57]: df[df['placement_exam_marks'] < lower_limit].shape
```

```
Out[57]: (0, 5)
```

```
In [58]: new_dff = df[df['placement_exam_marks'] < upper_limit]
```

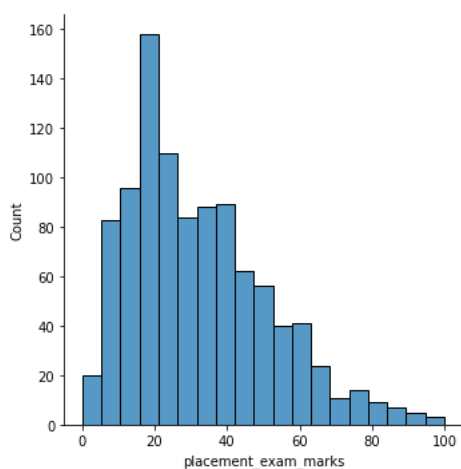
Trimming

```
In [59]: new_dff.shape
```

```
Out[59]: (985, 5)
```

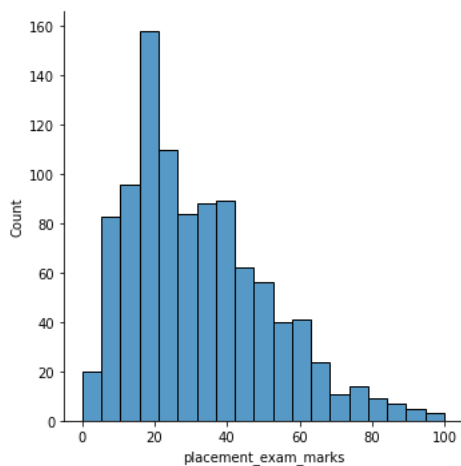
```
In [60]: sns.displot(df['placement_exam_marks'])
```

```
Out[60]: <seaborn.axisgrid.FacetGrid at 0x163a49d1970>
```



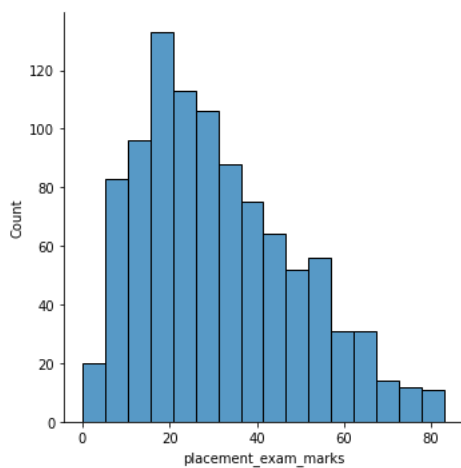
```
In [61]: sns.displot(df['placement_exam_marks'])
```

```
Out[61]: <seaborn.axisgrid.FacetGrid at 0x163a49d14c0>
```



```
In [62]: sns.displot(new_dff['placement_exam_marks'])
```

```
Out[62]: <seaborn.axisgrid.FacetGrid at 0x163a4b0c1f0>
```

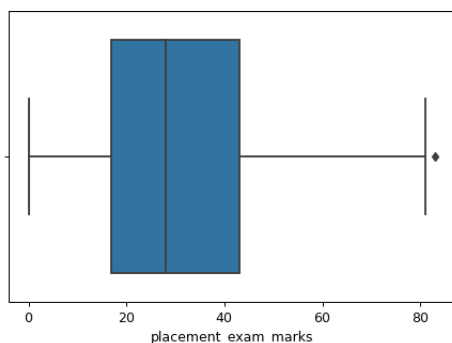


```
In [63]: %matplotlib notebook
```

```
In [64]: sns.boxplot(new_dff['placement_exam_marks'])
```

C:\Users\User38\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

warnings.warn(



```
Out[64]: <AxesSubplot:xlabel='placement_exam_marks'>
```

Capping

```
In [65]: new_dff_cap = df.copy()

new_dff_cap['placement_exam_marks'] = np.where(
    new_dff_cap['placement_exam_marks'] > upper_limit,
    upper_limit,
    np.where(
        new_dff_cap['placement_exam_marks'] < lower_limit,
        lower_limit,
        new_dff_cap['placement_exam_marks']

    )

)
```

```
In [66]: new_dff_cap
```

```
Out[66]:
```

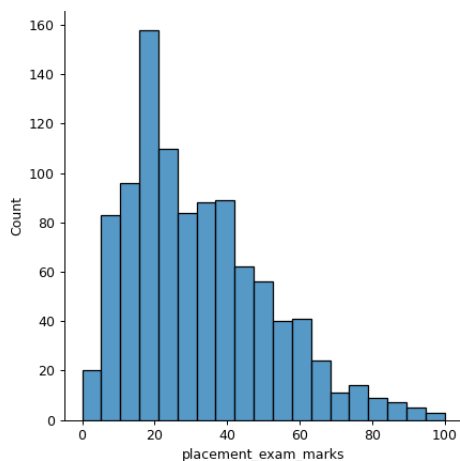
	cgpa	placement_exam_marks	placed	cgpa_score	cgpa_cap
0	7.19	26.0	1	0.371425	7.190000
1	7.46	38.0	1	0.809810	7.460000
2	7.54	40.0	1	0.939701	7.540000
3	6.42	8.0	1	-0.878782	6.420000
4	7.23	17.0	0	0.436371	7.230000
...
995	8.87	44.0	1	3.099150	8.808934
996	9.12	65.0	1	3.505062	8.808934
997	4.89	34.0	0	-3.362960	5.113546
998	8.62	46.0	1	2.693239	8.620000
999	4.90	10.0	1	-3.346724	5.113546

1000 rows × 5 columns

```
In [67]: new_dff_cap.shape
```

```
Out[67]: (1000, 5)
```

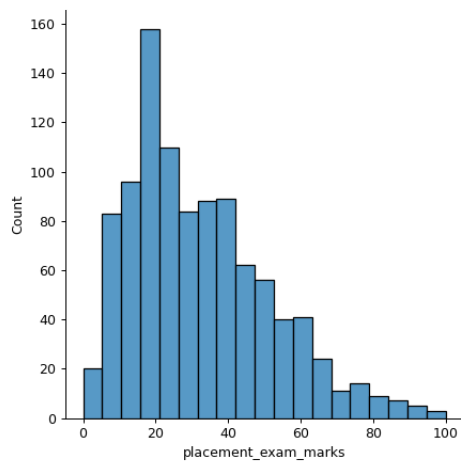
```
In [69]: sns.displot(df['placement_exam_marks'])
```



```
Out[69]: <seaborn.axisgrid.FacetGrid at 0x163a5cf1640>
```

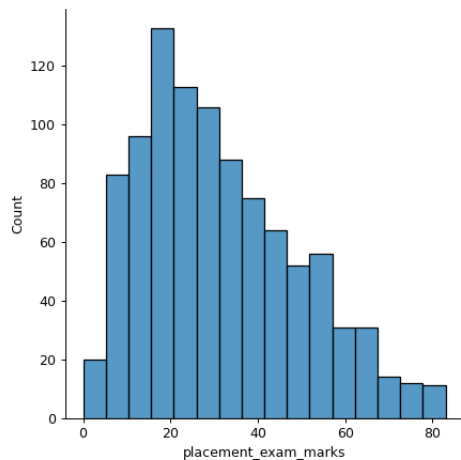


```
In [71]: sns.displot(df['placement_exam_marks'])
```



```
Out[71]: <seaborn.axisgrid.FacetGrid at 0x163a5d51d30>
```

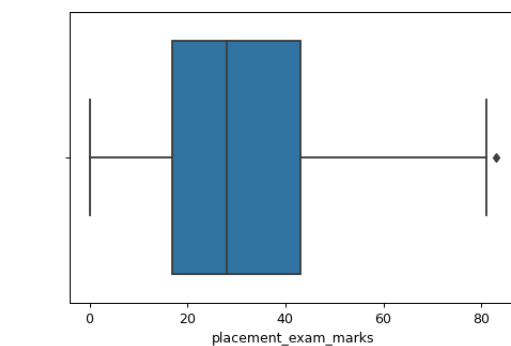
```
In [72]: sns.displot(new_dff['placement_exam_marks'])
```



```
Out[72]: <seaborn.axisgrid.FacetGrid at 0x163a5e66d30>
```

```
In [73]: sns.boxplot(new_dff['placement_exam_marks'])
```

C:\Users\User38\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.



```
Out[73]: <AxesSubplot:xlabel='placement_exam_marks'>
```

```
In [ ]:
```

